# Logistic Regression  (Week 3 lectures)

- classification

- $y$ is discrete.

## 1. Classification

ex.  email spam, online transactions, tumor.

· $y \in \{0, 1\}$ $\longrightarrow$ binary class.

negative class ↓

positive class ↓

· multi-class classification  $y \in \{0, 1, \ldots n\}$

· linear regression is not good for classification

1) although $y=0$ or $1$, $h_\theta(x)$ can output out of range, i.e $<0$ or $>1$.

$\Rightarrow$ logistic regression : output always in-range.  $0 \leq h_\theta(x) \leq 1$.
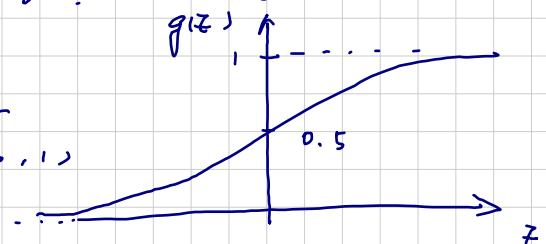
↓ classification.

## 2. Hypothesis Representation.

- want:  $0 \leq h_\theta(x) \leq 1$.    |  linear regression.  $\theta^T(x)$

$\Rightarrow h_\theta(x) = g(\theta^T x)$,  where $g(z) = \dfrac{1}{1+e^{-z}}$   (sigmoid function, logistic function )

↓ real number.

$= \dfrac{1}{1+e^{-\theta^T x}}$.

hence satisfies $h_\theta \in (0, 1)$

pick value for $\underline{\theta}$.

# Interpretation.

$h_\theta(x) =$ estimated probability of $y = 1$ on input $x$.

ex. $\underline{x} = \begin{pmatrix} 1 \\ tumor\ size \end{pmatrix}$ , $h_\theta(\underline{x}) = 0.7 \Rightarrow 0.7\%$ chance of being 1.

$\cdot\ h_\theta(x) = P(y = 1 \mid x; \theta)$

$\downarrow$ given x.
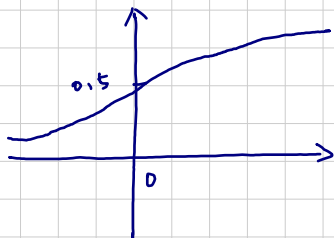
$\hookrightarrow$ parameterized by $\theta$.

$\cdot\quad y = 0$ or $1. \quad\Rightarrow\quad P(y=0 \mid x; \theta) + P(y=1 \mid x; \theta) = 1$.

## 3. Decision Boundary.

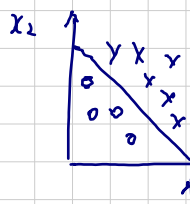$\cdot\quad h_\theta(x) = g(\theta^T x) ; \quad g(z) = \dfrac{1}{1 + e^{-z}}$

$= P(y = 1 \mid x; \theta)$

$\cdot\quad$ Suppose predict $\begin{cases} "y=1" & \text{if } h_\theta(x) \geq 0.5 &, \theta^T x \geq 0. \\ "y=0" & \text{if } h_\theta(x) < 0.5 &, \theta^T x < 0 \end{cases}$



$g(z) \geq 0.5$ when $z \geq 0$,

$h_\theta(x) = g(\theta^T x) \geq 0.5 \Longleftrightarrow \boxed{\theta^T x = 0}$

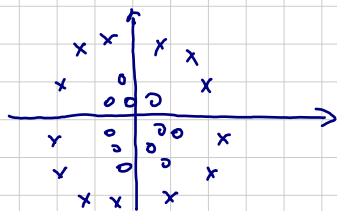$$h_\theta(x) = g(\theta^T \cdot x), \quad \text{say} \quad \theta = \begin{pmatrix} -3 \\ 1 \\ 1 \end{pmatrix}.$$

$$\Rightarrow y = 1 \quad \text{when} \quad \langle -3, 1, 1 \rangle \cdot \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} \geq 0.$$

$$\underbrace{\qquad\qquad}_{\text{equation of line}}$$

ex. suppose $\underline{\theta} = \begin{pmatrix} 5 \\ -1 \\ 0 \end{pmatrix}$ s.t. $h_\theta(x) = g(\underbrace{5 - x_1}_{\text{line}})$

$$5 - x_1 \geq 0 \Rightarrow x_1 = 5$$

Non-linear decision boundaries.

$$h_\theta(x) = (\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

suppose $\underline{\theta} = \begin{pmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{pmatrix}$

predict "$y = 1$" if $-1 + x_1^2 + x_2^2 \geq 0$.

$$\Rightarrow \text{equation of circle w/ radius } 1$$

# 4. Cost function.

- define optimization objective. (cost function)

- training set. $\{ (x^{(1)}, y^{(1)}), \quad (x^{(2)}, y^{(2)}), \quad \dots, \quad (x^{(m)}, y^{(m)}) \}$

- m examples $\quad x \in \begin{bmatrix} x_0 \\ \vdots \\ x_n \end{bmatrix}, \quad x_0 = 1, \quad y \in \{0, 1\}$

$\qquad \qquad \qquad \qquad e \mathbb{P}^{n+1}.$

- $h_\theta(x) = \dfrac{1}{1 + e^{-\theta^T x}}.$

$\qquad \qquad \qquad \hookrightarrow$ how to choose $\theta$ ?

- Recall:

linear regression. $\quad J(\theta) = \dfrac{1}{m} \cdot \sum\limits_{i=1}^{m} \dfrac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

$\qquad \qquad \qquad \qquad = \dfrac{1}{m} \sum\limits_{i=1}^{m} \text{cost} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

$\Rightarrow \text{cost} \left( h_\theta(x^{(i)}), y^{(i)} \right) = \dfrac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

$\Rightarrow \text{cost} \left( h_\theta(x), y \right) = \dfrac{1}{2} \left( h_\theta(x), y \right)^2$

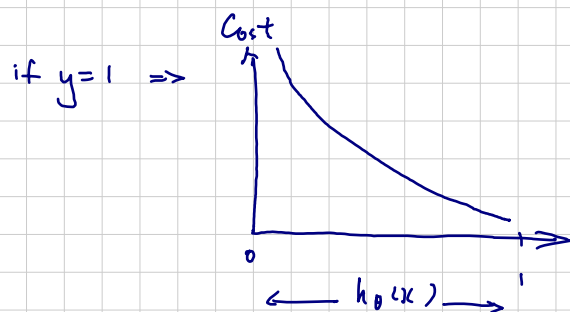$\qquad \qquad \qquad \hookrightarrow \dfrac{1}{1 + e^{-\theta^T x}}. \quad \Rightarrow$ non-convex for "square ost function"

- Logistic regression.



$\Rightarrow$ objective = create cost function s.t. $J(\theta)$ is convex.

$$\text{Cost}\left(h_\theta(x), y\right) = \begin{cases} -\log\left(h_\theta(x)\right), & y=1 \\ -\log\left(1-h_\theta(x)\right), & y=0 \end{cases}$$

if $y=1$ $\Rightarrow$



Cost

$\longleftarrow h_\theta(x) \longrightarrow$

- cost $= 0$  if $y=1$, $h=1$.

  But as $h_\theta(x) \to 0$, cost $\to \infty$.

- captures intuition if $h_\theta(x)=0$. predict.

  $P(y=1 \mid x; \theta) = 0$ but $y=1$

  $\hookrightarrow$ penalize learning algorithm w/ very

  large cost of $\infty$.

if $y=0$.



Cost

$h_\theta(x)$

$-\log(1-z)$

5. Simplified Cost function and gradient descent.

- $\text{Cost}(h_\theta(x), y) = -y \cdot \log(h_\theta(x)) - (1-y) \log(1 - h_\theta(x))$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} \sum_{i=1}^{m} y \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))$$

$\hookrightarrow$ maximum likelihood estimation.

To fit parameters $\theta$, $\min_\theta J(\theta)$

To make prediction given new $x$, output $h_\theta(x) = \dfrac{1}{1 + e^{-\theta^T x}}$ . $P(y=1 \mid x; \theta)$

- Gradient Descent.

$\min_\theta J$ : repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

}

$\dfrac{\partial}{\partial \theta} J(\theta) = \dfrac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$

$\downarrow$

def'n has changed

from linear regression!

$\cdot \quad \theta = \begin{pmatrix} \theta_0 \\ \vdots \\ \theta_n \end{pmatrix}$ , $x^{(i)} = \begin{pmatrix} x_0^{(i)} \\ x_j^{(i)} \\ \vdots \\ x_n^{(i)} \end{pmatrix}$ $\quad j \to$ feature #.

$\quad\quad\quad\quad\quad\quad$ one
$\quad\quad\quad\quad\quad\quad$ sample

- $\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$

$\underline{\theta} = \underline{\theta} - \frac{\alpha}{m} \cdot \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$

$\downarrow$

$(n+1) \times 1$.

# Advanced optimization

. cost function. $J(\theta)$, want $\min_{\theta} J(\theta)$.

. gradient descent          Repeat $\{$  $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$ $\}$

# Multiclass Classification.

ex. Email Foldering / Tagging    $\rightarrow$  $y = 1, 2, 3, 4$

work   friends

family   hobby.

## One-vs-all. (one-vs-rest)

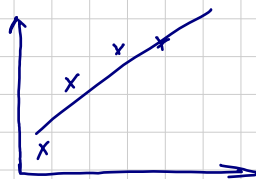- binary  (this group or not).
- iterate through every group.

. fit 3 classifiers:     $h_\theta^{(i)}(x) = P(y=i \mid x; \theta)$    $(i=1,2,3)$

$\hookrightarrow$ Train a logistic regression classifier. $h_\theta^{(i)}(x)$ for each class $i$ to predict.

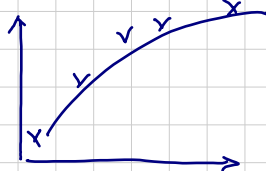Lecture 7.   Solving the Problem of overfitting.

1. The problem of overfitting.
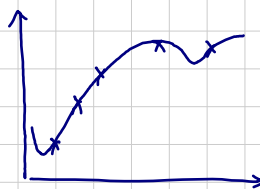
(1)   Underfit, high bias

$$\theta_0 + \theta_1 x$$

(2)   "Just right"

$$\theta_0 + \theta_1 x + \theta_2 x^2$$

(3)   Overfit, high variance

↳

tw fitting, can fit
any data

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

- two many features


Address overfitting

(1) Reduce number of features

- manually select which features to keep.
- model selection algorithm

(2)  Regularization

- keep all the features, but reduce magnitude / values of

parameters $\theta_j$.

- Works well when we have a lot of features, each of which
contributes a bit to predicting y.

# 2. Cost Function

· Intuition.

$$\theta_3 \approx 0, \quad \theta_4 \approx 0$$

suppose we penalize and make $\theta_3, \theta_4$ small.

$$\rightarrow \min_\theta \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + 1000\,\theta_3^2 + 1000\,\theta_4^2$$

· small values for some $\theta_i$ $\longrightarrow$ simpler hypothesis
$i = 0, \ldots, n$.

(small values
for all $\theta_i$)
$\longrightarrow$ less prone to overfitting.

· Example - housing

① features: $x_1, x_2 \ldots, x_{100}$

② Parameters: $\theta_0, \theta_1 \ldots \theta_{100}$.

$\rightarrow$ we don't know what parameter to "shrink"

$\Rightarrow$
shrink all parameters.

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

Ⓐ
we want to
Train data fit

$\lambda$ - regularization
parameter.
balances btwn
objectives Ⓐ and Ⓑ.

$\theta_0$ is discluded by
conventions

Ⓑ
we want to keep
$\theta$ small

· If $\lambda$ is set too large, all $\theta_i \approx 0$, then $h_\theta(x) \approx \theta_0$.

$\rightarrow$ At flat horizontal line

$\Rightarrow$ underfitting. (too high bias)

3. Regularized  Linear  Regression.

(1)

Gradient  Descent

repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{M} \left[ h_\theta(x^{(i)}) - y^{(i)} \right] x_0^{(i)}$$

$$\left[ -D \frac{\partial}{\partial \theta_0} J(\theta) \right.$$

$$\theta_j := \theta_j - \left\{ \alpha \frac{1}{m} \sum_{i=1}^{M} \left[ h_\theta(x^{(i)}) - y^{(i)} \right] x_j^{(i)} + \frac{\lambda}{m} \theta_j \right\}$$

$$j = 1, 2, 3 \dots h$$

$$\frac{\partial}{\partial \theta_j} J(\theta), \text{ regularized}$$

}

$$\Rightarrow \quad \theta_j := \theta_j \left( 1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^{M} \left[ h_\theta(x^{(i)}) - y^{(i)} \right] x_j^{(i)}$$

$$\underbrace{\qquad\qquad}$$

$$1 - \alpha \frac{\lambda}{m} < 1 .$$

$$\downarrow$$

small

$$\underbrace{\qquad\qquad\qquad}$$

original gradient descent

① make $\theta_j$ a bit smaller

② perform regular update.

(2) Normal Equation.

Design Matrix. $\underline{X} = \begin{pmatrix} -\ (x^{(1)})^T\ - \\ \vdots \\ -\ (x^{(m)})^T\ - \end{pmatrix}$ $\qquad \min_{\theta} J(\theta)$.

$\qquad\qquad\qquad\qquad m \times (n+1)$

$y = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{pmatrix}$

$\qquad\qquad\qquad (m \times 1)$.

$\dfrac{\partial}{\partial \theta_j} J(\theta) \overset{set}{=} 0 \qquad \Rightarrow \qquad \theta = \left( X^T X + \lambda \begin{pmatrix} 0 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix} \right)^{-1} X^T y$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \downarrow$

$\qquad\qquad\qquad\qquad\qquad\qquad e.g.\ n=2, \qquad (k+1) \times (k+1)$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$

Non-invertibility.

$\qquad\qquad$ # features.

Suppose $m \leq n$
$\qquad$ # examples

$\qquad\qquad\qquad\qquad\qquad \theta = (X^T X)^{-1} X^T y \rightarrow$ non-invertible / singular.

$\qquad\qquad\qquad\qquad\qquad \begin{cases} pinv. & pseudo\ inverse \\ \\ inv. & inverse \end{cases}$

If $\lambda > 0$,

$\qquad\qquad \theta = \left( \underbrace{ X^T X + \lambda \begin{pmatrix} 0 & & 0 \\ & 1 & \\ & & 1 \\ 0 & & 1 \end{pmatrix} }_{\text{invertible.}} \right)^{-1} X^T y.$

# 4. Regularized Logistic Regression.

$$J(\theta) = - \left( \frac{1}{m} \sum_{i=1}^{m} y^{(i)} \log \quad \cdots \quad + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2 \right)$$

Repeat {

$$\theta_0 := \theta_0 - \frac{\alpha}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \left( \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)} \right) + \frac{\lambda}{m} \theta_j$$

$$j = 1 \cdots n.$$

}

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

↳ regularized

$$h_\theta(x) = \frac{1}{1 + e^{\theta^T x}}.$$

# Advanced Optimization.

$$\underline{\theta} = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{pmatrix} \begin{matrix} - \text{theta}(1) \\ \vdots \\ - \text{theta}(n+1) \end{matrix}$$

function [jval, gradient]
= costfunction (theta).

↳ fminunc ↀ@ costfunction )

$$\frac{\partial}{\partial \theta_0} J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_0^{(i)}.$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \left[ \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)} \right] + \frac{\lambda}{m} \theta_j$$